

Lec 7

* clustering

↳ Problem of finding a hidden structure within unlabeled data.

الفكرة هنا انك بتقسم ال (data) لجزءات، كل جزء هو متشابهة لبعضها بتكون (cluster).

* K-means clustering

ex
↳ height, weight and average lifespan of animals.

↳ used for: clustering numerical data.

Input: numerical

↳ Euclidian distance: distance must be ~~existed~~ defined over variable space.

output: Centroid

المرکز دي (cluster)

~~↳ contents~~

use cases

→ introduction to classification (Discover classes)

→ exploratory technique

↳ Discover structure in data.

↳ Summarize properties of each cluster.

Algorithm

① له ادلة حاجة بتختار K بشكل عشوائي (3 نقط مثلاً)

② بعد كده بتسوي القيم للأقرب لكل نقطة منهم فيعملوا (cluster) مع بعض.

③ هتعمل إعادة حساب للـ (centroid) عشوائية تعدد واحد جديد.

له هتسب الـ (mean) لكل النقط الموجودة حوله الـ (cluster)

④ هتعيد الخطوة 2، 3 لحد تلاقي الـ (Centroids)

لا تتغير.

شرح الـ (algorithm) كدهور

موجود في صفحة 17 من 18
في معاينة

الـ (output) اللي بتطلع بيها

له الـ (center) بتاع آخر (cluster).

له ~~الـ (center)~~ آخر فرفه بتفرضه للـ (dataset) في آخر (cluster).

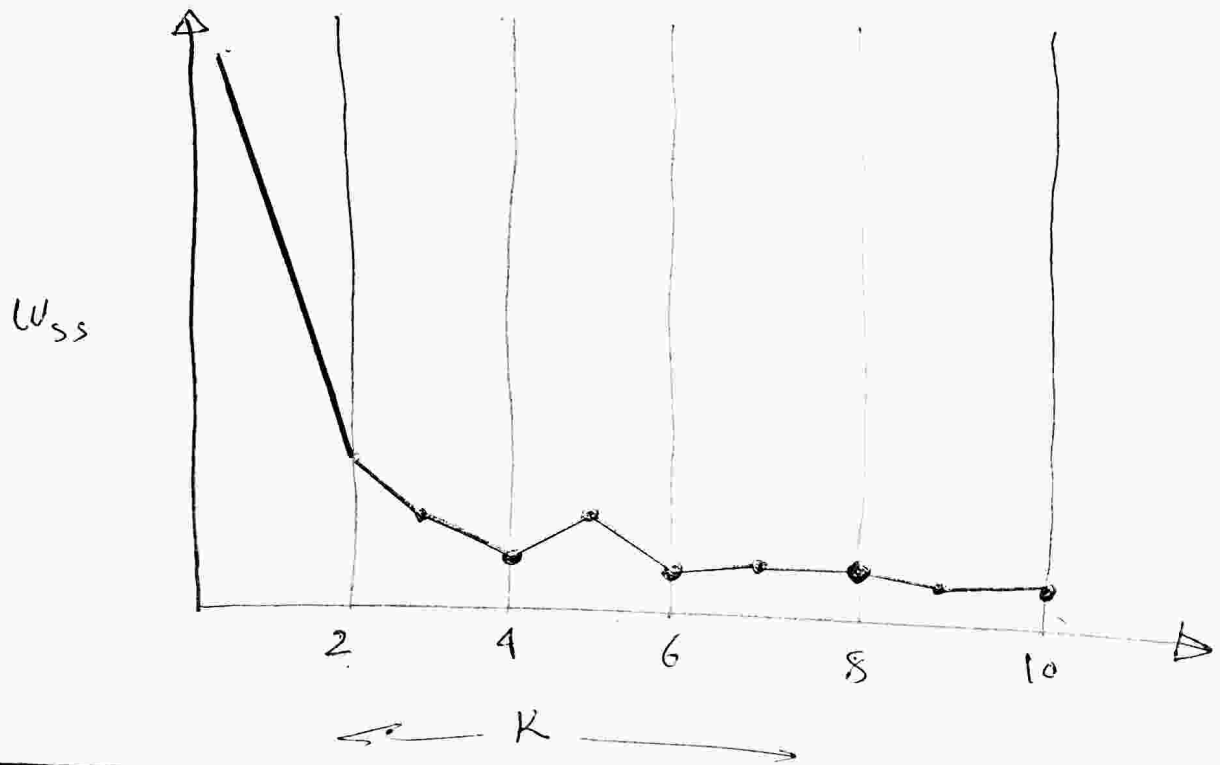
Picking K

Heuristic → Find the "elbow" of within-sum-of squares (WSS), plot it as function of K

$$WSS = \sum_{i=1}^K \sum_{j=1}^{n_i} |X_{ij} - C_i|^2$$

$K \rightarrow$ no. of clusters
 $n_i \rightarrow$ no. of points in i th cluster
 $C_i \rightarrow$ centroid of " "
 $X_{ij} \rightarrow j$ th point of " "

"Elbows" at $K = "2, 4, 6"$



مع كل تغير نحل N (evaluate) (model) فتكون:

(1) هل ال (clusters) باين انها منمضلة عالأقل في بعض الرسومات ولا لا ؟

(2) هل عندك (clusters) فيها (Points) قليلة -
لم جرب تقال ~~مع~~ قيمة K .

(3) هل فيه (Centroids) قريبة من بعضه
لم جرب تقال قيمة K .

K-means clustering

Reasons to choose	Cautions.
Easy to implement	Doesn't handle categorical variables.
• Easy to assign new data to existing clusters. ↳ which is the nearest cluster center	↳ sensitive to initialization (first guess)
↳ Concise output ↳ coordinates the K cluster centers.	• not scale invariant ↳ variables should all be measured on similar or compatible scales
	↳ Not always desirable ↳ tends to produce "round" equi-sized clusters.

Lec: 8

Association rules

↳ unsupervised learning method

↳ used to discover relationships within data.

قواعد الارتباط (Association rules) ←
في قواعد البيانات (databases)

→ Apriori → works on frequent itemset (set of items that appears together).

* Apriori Property

↳ Any subset of a frequent itemset is also frequent.

مثال لو عندك itemset C لو عندك L
"shoes, Purses" لو عندك L
وهي قولنا $\%50 = (\text{Support})$

لو عندك $\%50$ من (transactions) لو عندك (itemset) لو
فrequent itemset ← L

→ if 50% of itemsets have {shoes, Purses} in them then at least 50% of transactions with have either {shoes} or {Purses} in them → This is Apriori Property

Lift & Leverage

Lift

بموجب عدد المرات التي x, y سيظهروا فيها معاً أكثر من المتوقع لو أنها غير معتمدة على بعض.

→ هذه قياس لكيفية x, y ليها علاقة ببعضهم عن كونهم سيظهروا مع بعض.

$$\text{Lift}(x \rightarrow y) = \frac{\text{support}(x \cap y)}{\text{support}(x) * \text{support}(y)}$$

Leverage

الفرق بين احتمالية ظهور x, y معاً و احتمالية

إحد x, y غير معتمدة على بعض (وظهروا معاً بشكل متزامن)

$$\text{Leverage}(x \rightarrow y) = \text{support}(x \cap y) - \text{support}(x) * \text{support}(y)$$

Associative rules implementations

1) Market basket analysis

2) recommender systems.

3) discovering web usage patterns.

← مثال شرح ال (associations rules) من 11 إلى 17 في محاضرة

رقم ١ = فيه مسأله حسنه في ملف ثاني فيا بعد

من جهة حاجات بتقول بيها (check) و (data)

- 1) does data make sense?
- 2) make "test-set" from hold-out data.
- 3) Evaluate rules by lift or Leverage.

Notes

* support \rightarrow Percentage of transactions that contain L (set of items)

* Confidence \rightarrow Percentage of transactions that contain X, which also contain Y.

* output of apriori algorithm \rightarrow set of all rules $X \rightarrow Y$ with minimum support & confidence.

*Apriori

reasons to choose	Cautions
→ Easy to implement	*requires many database scans
→ uses clever observation to Prune Search space (Apriori Property)	*EXponential time complexity.
→ Easy to Parallelize	→ Addressed with Lift and Leverage measures
	to minimize

[8]